



COURSERA FINAL PROJECT

SPECIALIZED MODELS: TIME SERIES AND SURVIVAL ANALYSIS

Time Series Forecasting to predict grocery sales at Favorita stores

Project presented by

Raziel Amador Rios

to obtain the Coursera Certificate

Main object: produce a reliable forecast based on sales time-series data from a retail store (Favorita stores, an Ecuadorian-based company).

Barcelona, February 2023

Table of Contents

Table of Contents	iii
Time Series Forecasting	v
I Main objective	vi
II Data description	vii
III Data exploration and data cleaning	viii
III.1 Time series Decomposition	ix
III.2 Stationarity and seasonality analyses	x
IV Time-series Forecast	xi
IV.1 Prophet as our Forecast model	xii
IV.1.1 Prophet Hyperparameters	xiii
IV.2 Performance Metric	xiii
IV.3 Forecast results	xiii
Appendix	xv
V Supplementary information	xvi

Time Series Forecasting

I

Main objective

The **main object** of the project is to:

- Generate a reliable **time series forecast** based on sales from Favorita stores.

Favorita stores is a grocery retail store based in Ecuador. Producing an accurate forecast could lead to decreased food waste related to overstocking and improve customer satisfaction. All the code can be found in the following [GitHub repository](#), further information can be found in Supplementary information.

II

Data description

The dataset comes from **Kaggle**, named: [Store Sales - Time Series Forecasting](#). The Kaggle dataset contains sales data from *Corporación Favorita*, a large Ecuadorian-based grocery retailer. Kaggle provides you with 7 different files (see Table 1 for more details).

Number	File Name	Description
1	holiday_events.csv	Relevant holidays in Ecuador
2	oil.csv	Oil prices from 2013 to 2017
3	sample_submission.csv	submission example
4	stores.csv	Stores metadata
5	test.csv	Stores and family products
6	train.csv	Sales by store and product-family from 2013-01 to 2017-08
7	transactions.csv	Number of transactions by store

Table 1: Kaggle files description. Files are alphabetically sorted.

The training data represents 99% of the data, including dates from 2013-01-01 to 2017-08-16 (55.5 months), 54 stores placed in different cities within Ecuador, and 33 family-products (see Figure 1). The testing data includes dates from 2017-08-16 to 2017-08-31 (15 days).

Stores Summary					
54 Stores	5 Store types	17 Store clusters	33 Product families	16 States	55.5 Months

Figure 1: Summary of the training dataset. The cluster information denotes similarity between stores.

III

Data exploration and data cleaning

Performing an exploratory data analysis (eda) of the sales from the training dataset, we can observe that grocery I, beverages, and produce are the top 3 most consumed products (see Figure 2). Additionally, store type A, and cluster 5 are the most frequent among their classification (Figure 2).

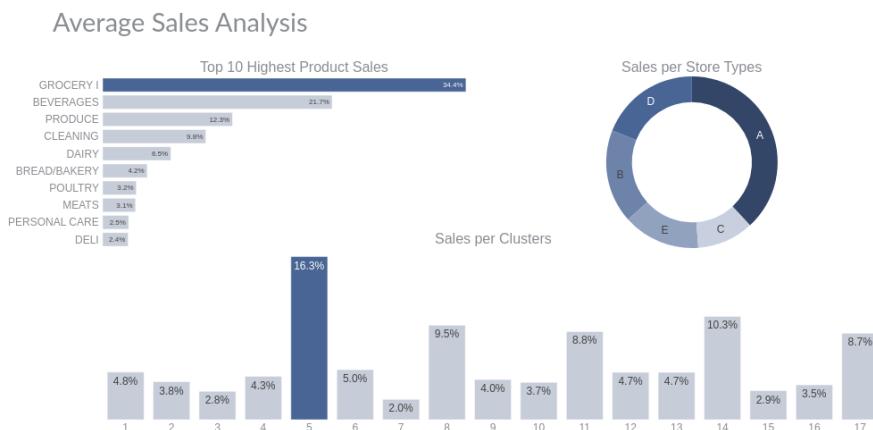


Figure 2: EDA of sales. The plot describes the sales by product, store type, and per cluster (left, right, and below, respectively). Darker blue represents higher sales.

The sales represented by Figure 3 shows low-peaks at the end of each year, which is explained because the stores are closed at New years. Moreover, we can observe a pattern at each year, with increased sales at the end of the year which overlaps with the Christmas eves, suggesting a seasonal pattern. These patterns are highlighted after smoothing the sales data with a 7 days moving average (the continuous red lines from the Figure 3).

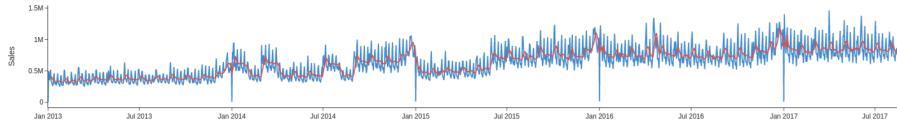


Figure 3: Target time series data. The y-axis represents the training sales from 2013-01-01 to 2017-08-15, and the x-axis shows the time in days (1,087 days). The sales were aggregated by all stores and products (see Figure 1). The raw data, and smoothed data (7 days moving average) are denoted by the continuous blue, and red lines, respectively.

III.1. Time series Decomposition

Time series are defined as a sequence of data organized in time order. The key components of time-series are: trend, seasonality, and residuals, displayed in Figure 4B, Figure 4C, and Figure 4D, respectively. Trend is the long-term direction of the time series, seasonality describes the periodic behavior such as holidays, and residuals represent the irregular fluctuations that we are not able to predict using the trend and/or seasonality.

In our target time series, we observe an upward trend, with a clear seasonality factor (this will be further analyzed below), and the residuals do not follow a random pattern. These insights will be used to select an adequate forecast model.

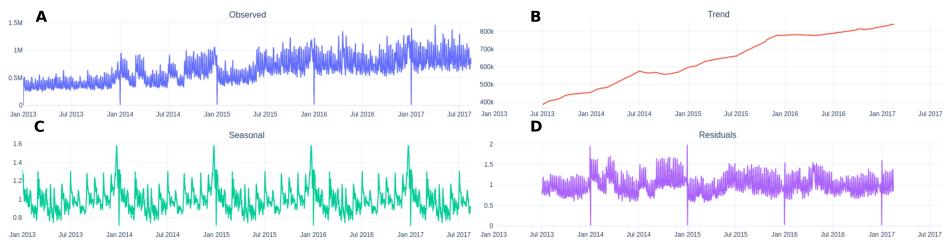


Figure 4: Time series decomposition. (A) Raw sales time series (sames as Figure 3). (B) Trend of sales. (C) Seasonality of sales. (D) Residuals of sales.

III.2. Stationarity and seasonality analyses

The stationary and seasonality are relevant components to select the adequate machine learning model (such as ARIMA, SARIMA, etc.) to generate reliable forecasting. A stationary time series is defined, if their statistical properties such as mean, and variance are all constant, and independent of time.

In consequence, we implemented a *Dickey-Fuller test* to assess stationarity. We obtained a *p-value* of 0.09, using a significance level of 0.05, we can reject the null hypothesis (the series is stationary) concluding that our series is non-stationarity.

In terms of seasonality, we aggregate the sales for each month and we can clearly see on average an increase of sales during December (Figure 5), which makes sense because Christmas and New years represent a season of the year with more transactions in other retail businesses. Furthermore, we highlighted February as the lowest sales month.

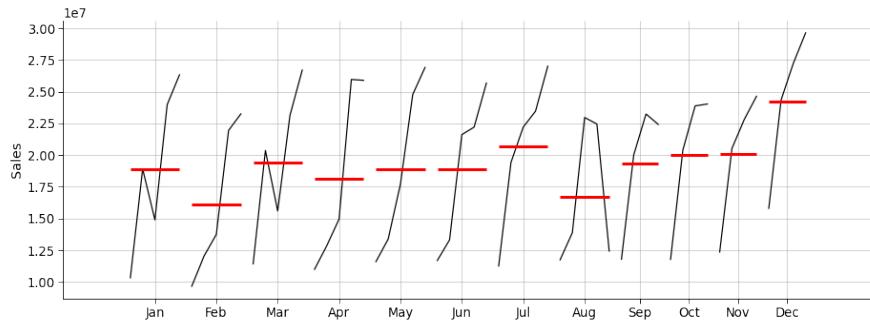


Figure 5: Analysis of seasonality. The y-axis represents the sales aggregated by month, x-axis shows the analyzed months (from January to December), the horizontal red lines denote the sales average by month (55.5 months).

Holidays are one of the most important factors for seasonality, and the Kaggle dataset provided us with relevant holidays and special events in Ecuador (see Table 1). Thus, we analyzed the effects of holidays on sales (see Figure 6). We observed a reasonable correlation between holidays and sales. Consequently, holidays must be incorporated on the forecast.

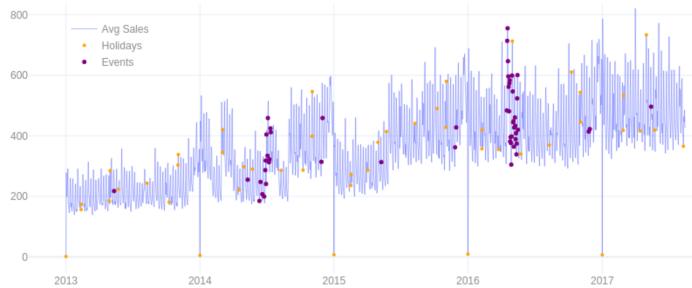


Figure 6: Holidays effect. Holidays and special events are represented as yellow, and purple dots, respectively.

IV

Time-series Forecast

The main goal of the project is to generate a reliable forecast using the sales data from Favorita stores. After the eda, we concluded that our data is not stationary, possesses a strong seasonal effect, and holidays is a factor that should be considered in the forecast. Therefore, the following models could be implemented:

- **SARIMAX:** Seasonal ARIMA using eXternal data. Further, ARIMA stands for Auto-Regressive (p: dependence on past values), Integrated (d: differencing) Moving Average (q: dependence on past forecast errors).
- **LSTM:** Long-Short Term Memory.
- **GRU:** Gated Recurrent Unit.
- **Facebook Prophet:** is a forecasting procedure based on a general additive model where non-linear trends are fit with yearly, weekly, and daily seasonality, plus holiday effects.
- **NeuralProphet:** is a hybrid forecasting framework based on PyTorch and trained with standard deep learning methods.

- **LightGBM** (or other extreme gradient boosted method, e.g. XGBoost): Light Gradient-Boosting Machine.

Ideally, we would like to implement all of the machine learning models mentioned above, do a benchmark taking into account: performance, computational-resources, training-time, and model-explainability; and compare their results. Nonetheless, due to time-constraints restrictions, we selected **Facebook Prophet** as our unique machine learning model to forecast the sales from Favorita stores.

The reasons to select Facebook Prophet include: a simple and flexible model (further described in below) that takes into account multiple human scale seasonality, holidays that occur at irregular intervals, trends that are non-linear, short training time, low computation-resource demands, a Python module to easily implement the model (under the hood uses Stan), and previous experience with the model.

IV.1. Prophet as our Forecast model

Facebook Prophet or Prophet (we are going to use Prophet from now on) is an automated forecasting procedure based on a general additive model where non-linear trends are fit with yearly, weekly, and daily seasonality, plus holiday effects.

According to Taylor *et al.*, Prophet works best with time series that have strong seasonal effects (which is our case see Figure 4C, and Figure 5), holidays effect (Figure 6), and typically handles outliers well. The additive models contains three main components: **1)** trend, **2)** seasonality, and **3)** holidays, which are combined in the following equation:

$$y(t) = g(t) + s(t) + h(t) + \varepsilon_t$$

Where $g(t)$ is the trend function which models non-periodic changes in the value of the time series, $s(t)$ represents seasonality (e.g., weekly and yearly seasonality), and $h(t)$ represents the effects of holidays which occur on potentially irregular schedules over one or more days. The error term ε_t represents changes not explained by the model; under the assumption that ε_t is normally distributed.

IV.1.1. Prophet Hyperparameters

Prophet possesses 5 hyperparameters which can be tuned through cross validation. The hyperparameters are: **1) changepoint prior scale**, **2) seasonality prior scale**, **3) holidays prior scale**, **4) seasonality model**, and **5) changepoint range** (this one not recommended to be tuned according to the Prophet developers).

In our work, only the *seasonality mode* (either additive or multiplicative) was tuned using cross validation with an horizon of 1 year (365 years), period of 180 days, and 730 days for training. Information about cross validation can be found on [the Prophet site](#) under the diagnostics section, or at Taylor *et al.* work.

IV.2. Performance Metric

To assess the performance of our models, we selected the mean absolute percentage error (MAPE) which is defined below:

$$MAPE = \frac{100\%}{n} \sum_{t=1}^n \left| \frac{A_t - F_t}{A_t} \right|$$

Where A_t is the actual value (true value) and F_t is the forecast result (predicted value). MAPE was selected as a performance metric because of its simplicity to represent the forecast errors within a percentage scale, where the lower the better.

IV.3. Forecast results

Appendix

V

Supplementary information

This work was written with emacs¹ using L^AT_EX², using only **Free and Open Source software**. All the computational analysis were carried out using Linux-based distributions. The figures were generated with Python (matplotlib³/seaborn⁴/plotly⁵) and Inkscape⁶.

Contact: razielar@gmail.com

¹<https://www.gnu.org/software/emacs/>

²<https://www.latex-project.org/>

³<https://matplotlib.org/>

⁴<https://seaborn.pydata.org/>

⁵<https://plotly.com/>

⁶<https://inkscape.org/>